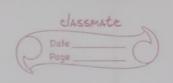# Practical - 3

Write code to simulate requests coming from clients and distribute them among the servers using load balancing algorithms.

→ what is load balancing?
- Load balancing is the process of distributing incoming network traffic across multiple servers.
- Purpose:-
  - No single servers get overloaded.
  - All servers are used efficiently.
  - Systems stay fast, reliable and scalable

→ Load balancing algorithms:

1) Round Robin :- Each server gets a turn in strict or order.
   - Maintains a list of available server.
   - Assign a new request to next server in the list and move that server to the end of the list.
   - This ensures that each server receives requests in a round - robin fashion.
   - Basic web applications.
   - Drawback: even if server is already busy, it still gets new request.

2) **least connections:**
- Server with the fewest active client gets ~~not~~ next request.
- Maintain a counter for the number of active connections on each server.
- This distributes the ~~workload~~ requests based on current work loads.
- Video streaming services, chat server where processing time can vary a lot.
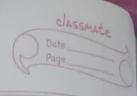- Drawbacks: Can cause overhead : needs to continuously monitor active connections for all servers. Slower if connections are constantly switching.

3) **Weighted round robin:**
- Server with higher weight gets more request.
- Assign weights to each server based on processing power, memory or other relevant factors.
- In a cloud environment with mixed server size, if some server ~~has~~ are more powerful (more RAM, CPU) they are assigned more work.
- Draw back: Fixed weight can't handle changing server performance.

→ **Real life applications:-**
- Google, face book, Youtube : have millions of users, at one requests distributed across thousands

of servers.
- Banking applications :- checking balance, removing
- Cloud Systems (AWS, Azure). / cash etc.
- Gaming servers (PUBG, fortnite).

→ Components in code :
1) ClientRequest class :-
   - Simulates a user request comming
     to the system.
   - Each request has :
     - request - id
     - arrival time
     - processing time.
2) Server class :-
   - Simulate a server that can handle
     client requests.
   - Each server keeps track of
     - Active connections → how many clients it
       is currently serving.
     - Total requests handled.
     - Weight.
   - When a server handles a request it.
     - Increments its active connections.
     - Sleeps (wait) for the processing
       time.
     - Decrements active connections
       after completing.
     - Tracking of active
       connection is not needed but
       its used for future upgrade
       or changes to last connections.

- It's harmless overhead - helps in statistics, monitoring.

3) Class LoadBalancer

☆ We are using weighted round robin.

→ Multithreading:-

→ Multithreading:- Each request is handled on a new thread → Simullating multiple users interacting with the system at the same time.

→ Locking (self. lock) : To prevent race conditions when multiple thread try to change data at the same time.

→ Randomness: Client arival times and processing times are random. ⊛